

CASE STUDY

Building a Prompt Engineering Toolkit

Uber



Introduction

- Crafting effective prompts is key for generating accurate outputs from LLMs (Large Language Models).
- To streamline LLM experimentation, Uber developed a centralized prompt engineering toolkit for building, managing, and executing prompt templates across models.
- The centralized toolkit helps create well-crafted prompt templates, enhanced with context from RAG (Retrieval-Augmented Generation) and runtime feature datasets.
- The toolkit enables prompt creation with system instructions, dynamic contextualization, large-scale batch offline generation (LLM inference), and prompt response evaluation.
- It also supports version control, collaboration, and safety measures, including hallucination checks, a standardized evaluation framework, and a safety policy to ensure responsible AI usage.
- Uber's case study provides an overview of the prompt template lifecycle, toolkit architecture, and its use in production at Uber.



Goals

Uber designed the prompt engineering toolkit to enable users to:

- Explore available LLM models.
- Search and discover prompt templates.
- Create and update prompt templates and tune parameters with access controls.
- Iterate on prompt templates through code reviews and version control.
- Generate batch offline generation by applying RAG and testing datasets with different LLM models for LLM inference.
- Evaluate prompts templates via offline evaluation through testing dataset and online evaluation.
- Monitor prompt template performance with LLMs in production to detect regressions and continuously monitor LLM quality through systematic tracking.



Prompt Engineering Lifecycle

The case study outlines two stages of the prompt engineering lifecycle: development and productionization, with users engaging with the prompt toolkit at each stage.

Development Stage

This stage includes three phases:

- 1) **LLM Exploration:** Users explore LLMs through a model catalog and Uber's GenAI Playground. The catalog provides descriptions, metrics, and usage guides for available models, while the playground allows users to test LLM capabilities and their applicability to specific use cases.
- 2) **Prompt Template Iteration:** Users define business needs, gather sample data, and create, analyze, and refine prompts. Auto-prompting makes suggestions so users don't need to create prompt templates from scratch, while a catalog of existing templates supports reuse. Iterative testing and revisions ensure alignment with objectives.
- 3) **Evaluation:** Users evaluate prompts against more extensive datasets to assess performance, leveraging either LLM-based evaluation (LLM as a judge) or custom code-based evaluators to ensure prompt effectiveness.

Productionization Stage

In this stage, only prompt templates that meet the evaluation threshold are productionized. Users monitor usage in production, collecting data on system usage to identify opportunities for improvements.

Architecture

- The prompt engineering toolkit integrates various components to support LLM deployment, prompt evaluation, and batch inference responses.
- It includes a Prompt Template UI/SDK for managing templates and revisions, connected to GetAPI and ExecuteAPI for interfacing with LLMs in the Model Catalog.
- Templates and models are stored in ETCD and UCS (object configuration storage) and utilized in the Offline Generation and Evaluation Pipelines.

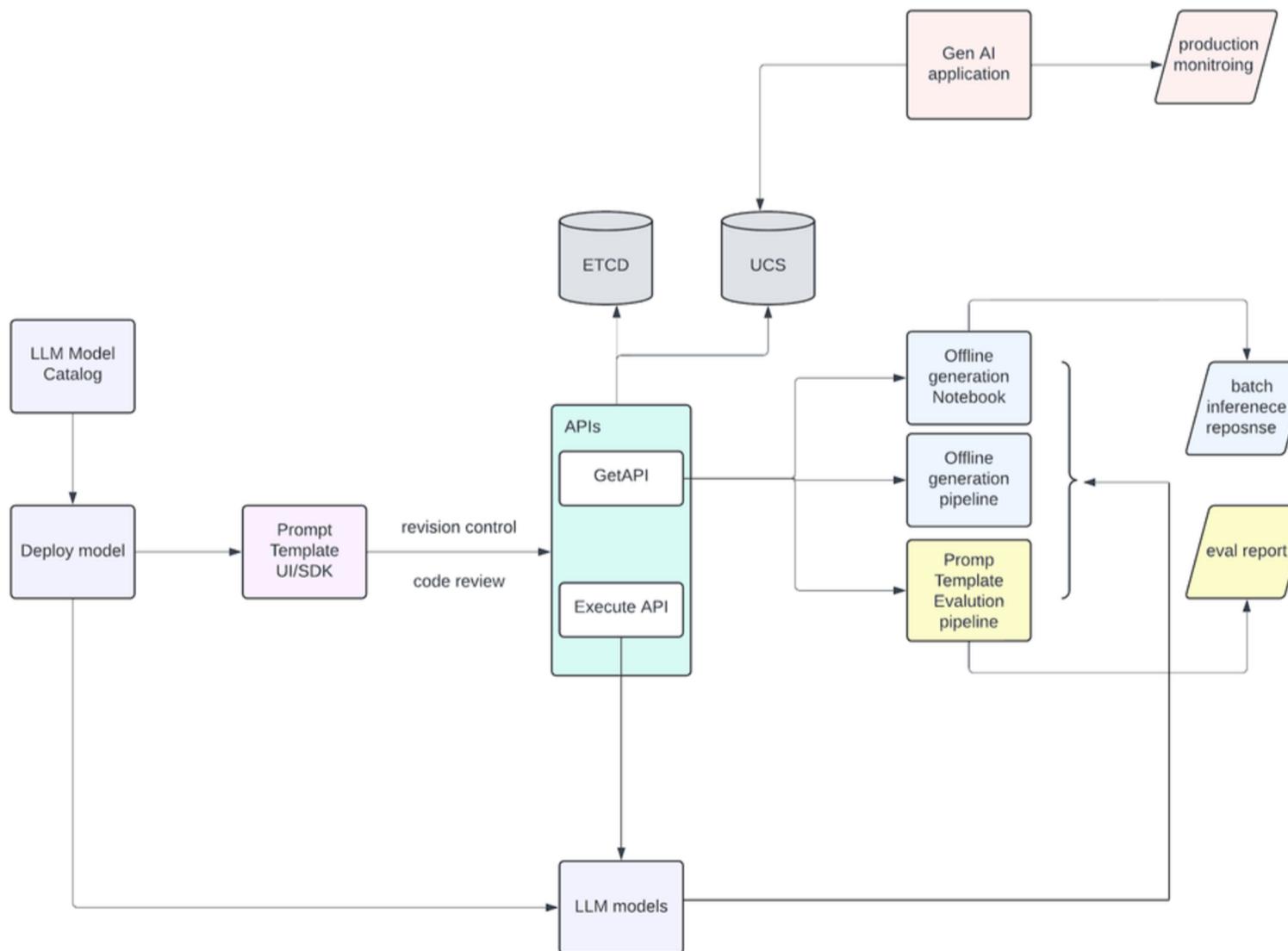


Figure 1: Overview of the prompt engineering toolkit.



UI/Design

- Users interact with the prompt toolkit in GenAI Playground, which integrates with the Model Catalog: a repository of Uber's LLMs with specifications, typical use cases, costs, and performance metrics.
- In the playground, users can select models, craft custom prompts, and adjust parameters (e.g. temperature) to evaluate responses during the ideation phase.

Model catalog

A list of available foundational and fine tuned models. [View](#)

Models

Search... [+ Add filter](#)

Azure OpenAI Service GPT-3.5 Turbo	Azure OpenAI Service GPT-3.5 Turbo 16K	Azure OpenAI Service GPT-4	Azure OpenAI Service GPT-4 32K
Foundation Language GPT-3.5, developed by OpenAI, is a part of the GPT series of language models. It's an evolution of the GPT-3 model, optimized for chat and other natural language processing tasks. The GPT-3.5 series includes models like text-davinci-002, text-davinci-003, and the most capable, gpt-35-turbo. These models have been trained using various strategies...	Foundation Language GPT-3.5, developed by OpenAI, is a part of the GPT series of language models. It's an evolution of the GPT-3 model, optimized for chat and other natural language processing tasks. The GPT-3.5 series includes models like text-davinci-002, text-davinci-003, and the most capable, gpt-35-turbo. These models have been trained using various strategies...	Foundation Language GPT-4 is a large-scale, multimodal model with the ability to accept image and text inputs, generating text outputs. It demonstrates human-level performance across various professional and academic benchmarks. The model employs a Transformer-based architecture, pre-trained on diverse data to predict the next token in a document, showcasing...	Foundation Language GPT-4 is a large-scale, multimodal model with the ability to accept image and text inputs, generating text outputs. It demonstrates human-level performance across various professional and academic benchmarks. The model employs a Transformer-based architecture, pre-trained on diverse data to predict the next token in a document, showcasing...
See Details >	See Details >	See Details >	See Details >
Azure OpenAI Service GPT-4 Turbo	Facebook Research Llama 2	Google AI PaLM 2 for Chat (Chat-Bison)	
Foundation Language GPT-4, developed by OpenAI, represents the latest advancement in the GPT series of language models, building upon the success of GPT-3 and GPT-3.5 with significant improvements in understanding and generation capabilities across a wide range of natural language processing tasks. This model iteration includes advanced models optimized for...	Foundation Language Llama 2 is a state-of-the-art, large-scale language model with 70 billion parameters. It's designed to understand and generate human-like text based on the input provided. The model is a Transformer-based architecture and is a successor to the original Llama model, offering enhanced performance and capabilities.	Foundation Language PaLM 2 for chat (Chat-Bison) is a state-of-the-art, large-scale language model optimized for natural language tasks like chatbots or help agents. It excels at language understanding, generation, and multi-turn conversations. **Cost** For general information regarding the costs associated with this Model, please refer to the publicly available pricing details provided...	
See Details >	See Details >	See Details >	
Google AI Text Embeddings Gecko	Mixtral 8x7B	OpenAI GPT-3.5 Turbo	
Foundation Embedding textembedding-gecko is a model that supports text embeddings, a NLP technique that converts textual data into numerical vectors. These vectors can be processed by machine learning algorithms, especially large models. The embeddings capture the semantic meaning and context of the words they represent. There are multiple versions of this...	Foundation Language Mixtral 8x7B Instruct is a Sparse Mixture of Experts (SMoE) model fine-tuned to follow instructions. Mixtral has the same architecture as Mixtral 7B, with the difference that each layer is composed of 8 feedforward blocks (i.e. experts). For every token, at each layer, a router network selects two experts to process the current state and combine their outputs. Even...	Foundation Language GPT-3.5, developed by OpenAI, is a part of the GPT series of language models. It's an evolution of the GPT-3 model, optimized for chat and other natural language processing tasks. The GPT-3.5 series includes models like text-davinci-002, text-davinci-003, and the most capable, gpt-35-turbo. These models have been trained using various strategies...	
See Details >	See Details >	See Details >	

Explore Gen AI

Explore the capabilities of different LLMs available at Uber, use case applicability and fine-tune the interactions with the models using prompts to achieve desired output. [View](#)

System instructions [+ Use prompt builder](#)

Act as an AI go code assistant

Configuration

Model [GPT-4o](#)

OpenAI is now the recommended production provider on Michelangelo Gen AI Gateway. Refer to EngSec guidelines for PII reduction requirements for using L1/L2 data. [View](#)

Task Chat Completion

Temperature 0.2 **Max Tokens** 1024

Top-K 40 **Top-P** 0.8

Parameters [+ Add more](#)

Responses

2024-03-09 12:37:45 PM [Show Input](#) [Reuse Settings](#)

Writing unit tests in Go is straightforward thanks to its built-in testing package. Here's a step-by-step guide on how to write a unit test in Go:

[Create a Test File](#)

Your test file should be in the same package as the code you're testing.

Figure 2: The model catalog UI.

Figure 3: The prompt template edit page.



Prompt Template Creation

The prompt engineering toolkit includes an auto-prompt builder that helps users create prompts and discover techniques for their specific use cases. Built on Uber's Langfx framework (an Uber internal Langfx service built on LangChain), it follows these steps:

1. **Best Practices:** Integrates prompt engineering best practices for context RAG retrievers.
2. **Templates and Examples:** Offers templates with detailed instructions and few-shot examples to help LLMs craft user prompts.
3. **LLM-Powered Generation:** Leverages LLMs to auto-generate and suggest tailored prompts.

The prompt builder leverages techniques such as chain of thought (CoT), automatic Chain of Thought, prompt chaining, tree of thought, Automatic prompt engineer, multimodal CoT prompting (see full case study for details on each)

Prompt Template Revision Control

- Prompt template iteration follows review practices. Users modify instructions and model parameters, testing responses with a dataset. Each iteration of a prompt template requires a review, and once approved, a new revision is created.
- To prevent accidental changes in production, users can deploy prompt templates with a specific revision tag. The deployed templates are managed through ObjectConfig, Uber's internal configuration system, and are automatically fetched by the production service upon deployment.

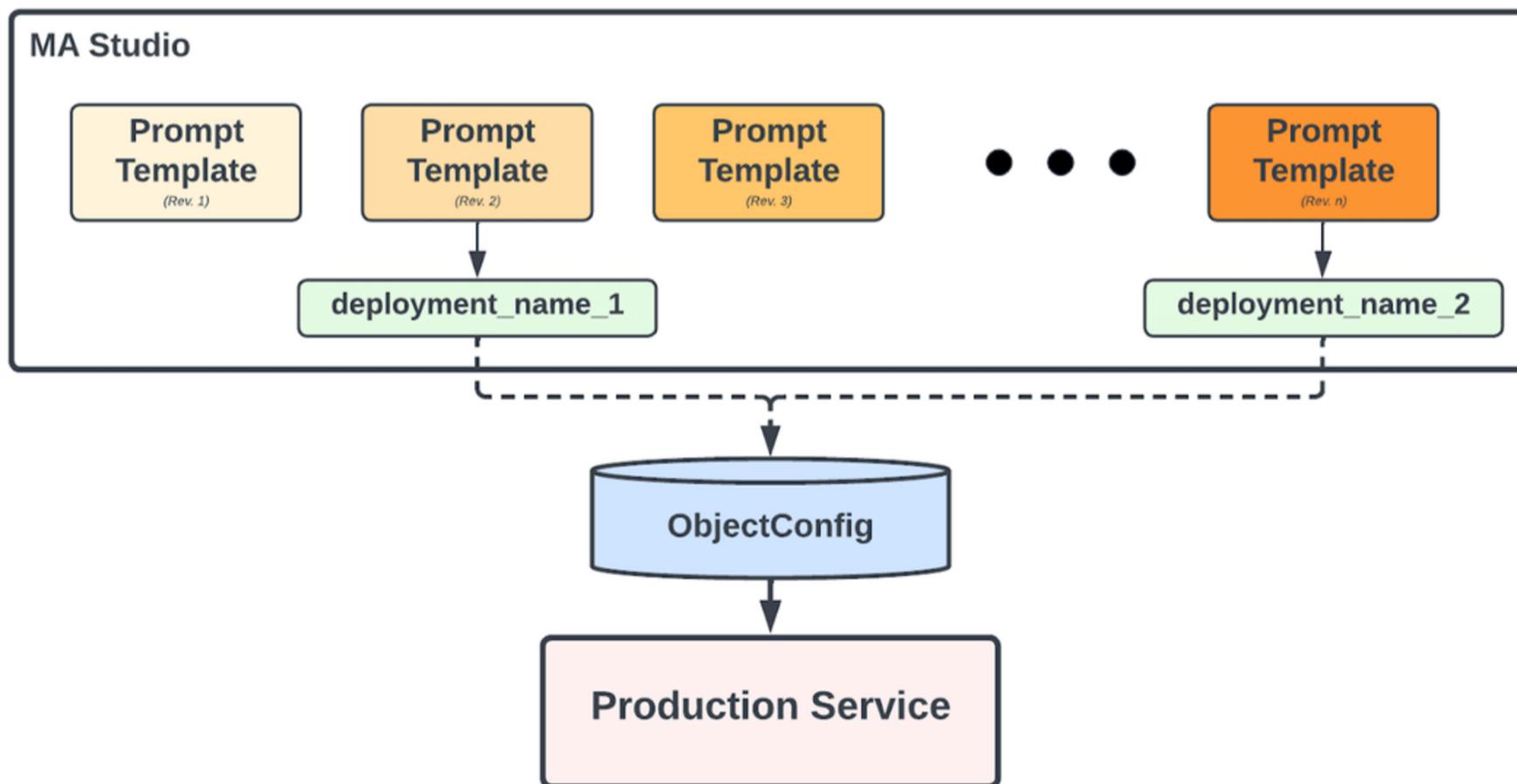


Figure 4: Prompt template revision control.

Prompt Template Evaluation

For evaluating prompt templates, two methods are employed by Uber:

1. **LLM as Judge:** Useful for tasks requiring subjective quality, such as generating text that should be engaging, persuasive or stylistically specific.
2. **Custom, user-defined code:** Specific metrics and criteria are coded to evaluate LLM performance, offering particularly tailored evaluation if required.

The toolkit provides evaluation prompt templates, each with instructions, examples, metrics and the format for responses. Templates are evaluated using either a golden dataset with labeled data or a dataset derived from production traffic. Aggregated metrics from large datasets offer a high-level comparison of template effectiveness.

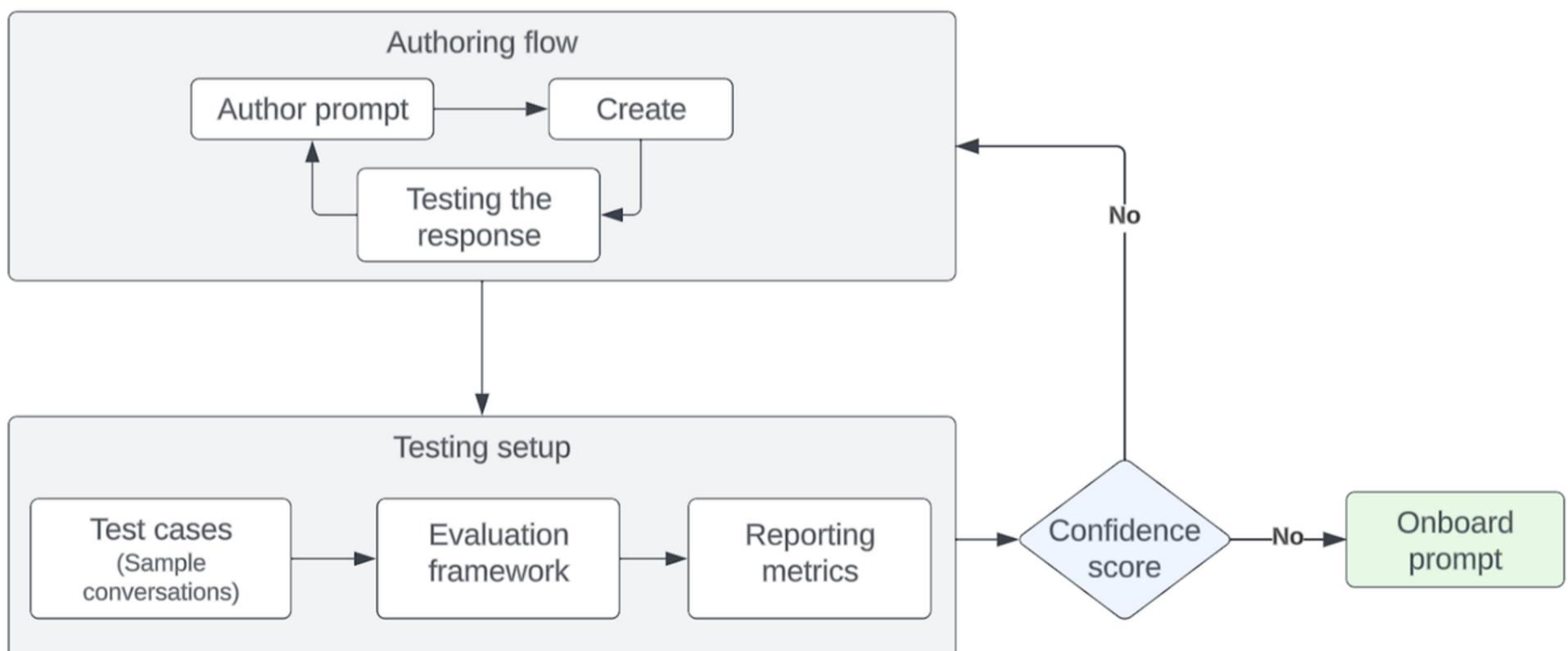


Figure 5: Prompt evaluation flow.

Production Use Cases at Uber Offline LLM Service

- The LLM Batch Offline Generation pipeline enables large-scale LLM inference for tasks like rider name validation, which checks the legitimacy of usernames. This pipeline processes existing and new usernames asynchronously, generating responses in batches.
- In MA Studio (Uber's AI platform), users set up the pipeline by selecting a relevant dataset, which is then dynamically applied to the prompt template. For example, the template "Is this {{user_name}} a valid human name?" is populated with usernames from the dataset, generating custom prompts for each entry.

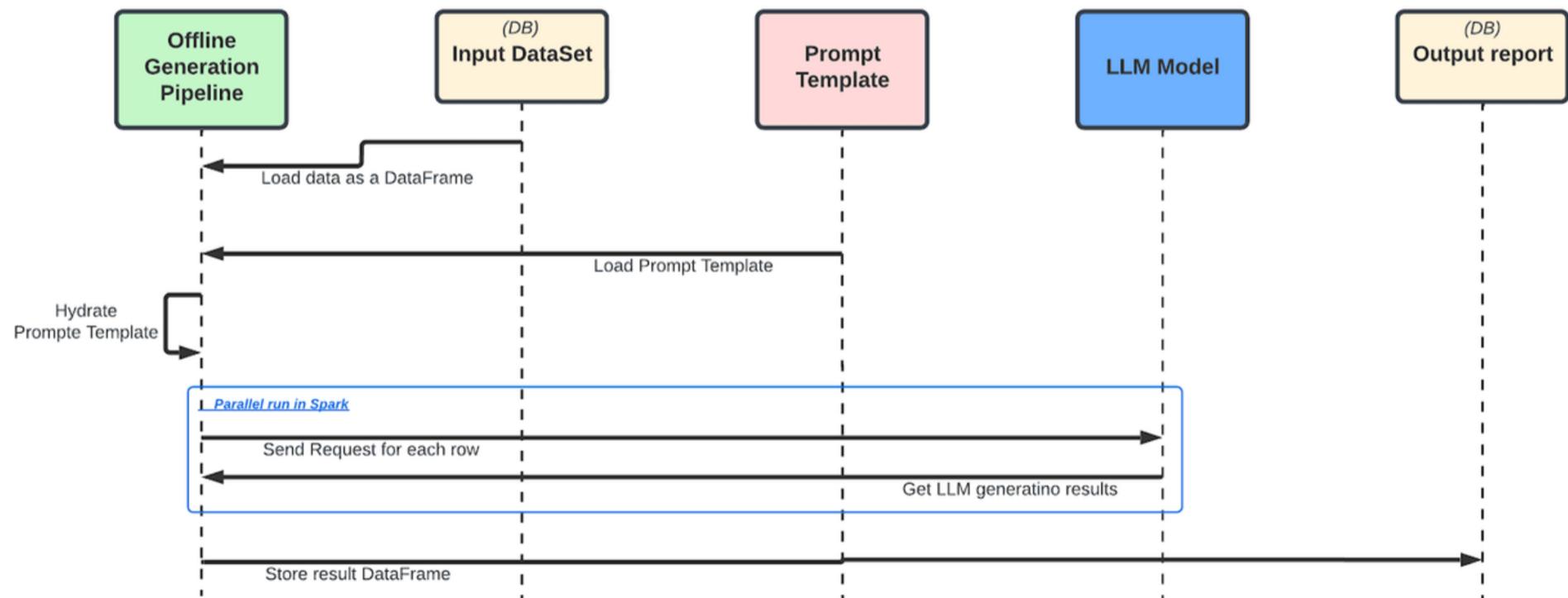


Figure 6: Prompt offline generation.

Production Use Cases at Uber Online LLM Service

The prompt template contains dynamic placeholders that are substituted with runtime values provided by the caller. Currently, only string-type substitution is supported using Jinja-based template syntax.

The service also supports fan-out capabilities across prompts, templates, and models:

1. **Templates:** The API formats payloads into vendor-specific structures from the generic data model. Multiple templates, such as chat and text completion, are supported.
2. **Prompts and Models:** Prompts are linked to specific models and templates. The service fetches and executes the required parameters via genAI APIs.

For example, in a summarization use case, when a support ticket is handed off to a new agent, the service generates a summary to provide context, eliminating the need for the new agent to ask the customer to repeat the issue.

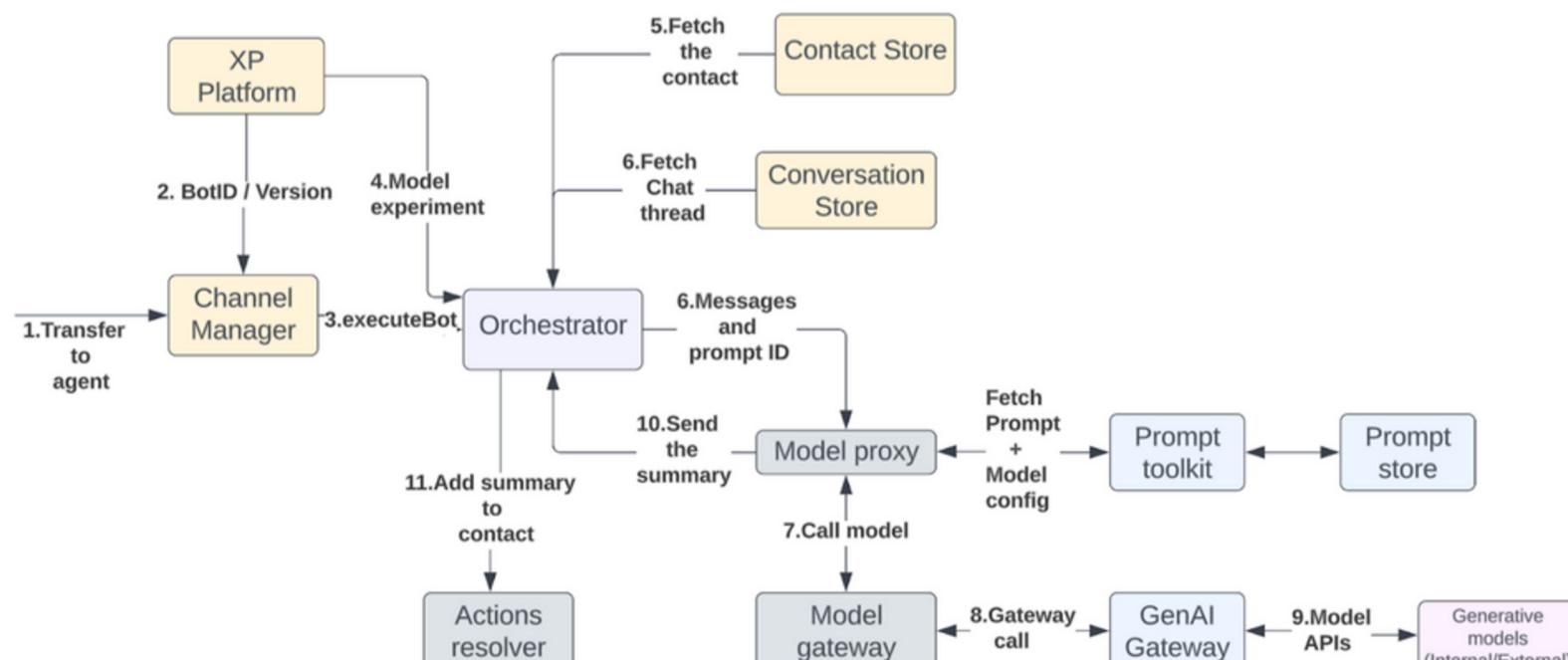


Figure 7: Prompt execution and summarization flow.



Conclusion

- The prompt engineering toolkit provides a framework to enhance LLM interactions at Uber, supporting development and production stages.
- From exploring LLM capabilities in the GenAI Playground to creating and iterating prompt templates, it enables users to leverage LLMs effectively.
- The toolkit's architecture ensures a systematic approach to prompt design, combining advanced guidance and evaluation methods for quality results.
- The prompt template lifecycle, from development to production and monitoring, enables rigorous testing and performance optimization.

Attribution

This document summarizes key points from Uber's case study, "Introducing the Prompt Engineering Toolkit" (November 2024). Read the full case study here:
<https://www.uber.com/en-CA/blog/introducing-the-prompt-engineering-toolkit/>

Want to level up with Generative AI? Follow



Lewis Walker



Repost this to help your network